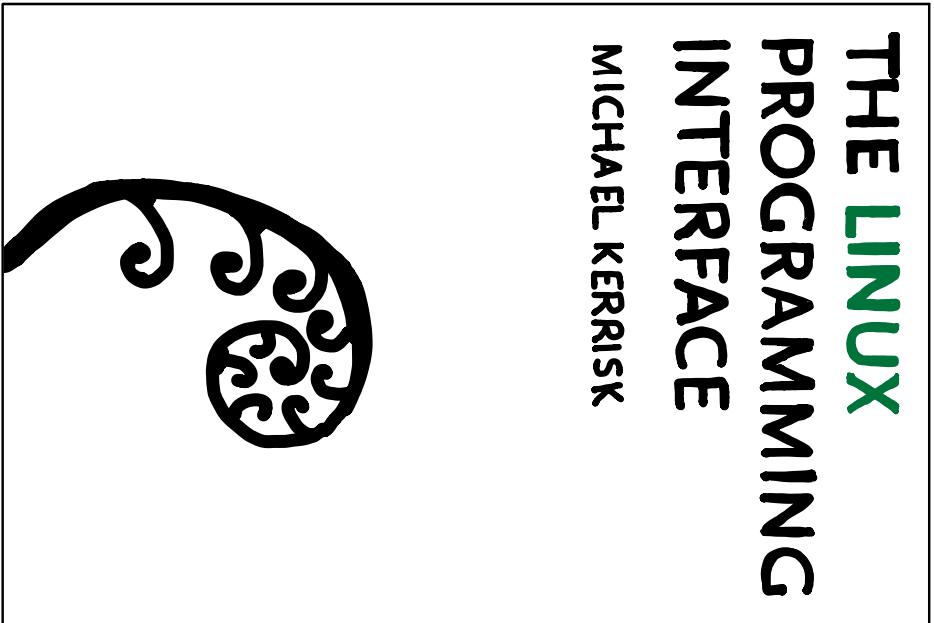


love this?

find more awesome zines at
→ wizardzines.com ←



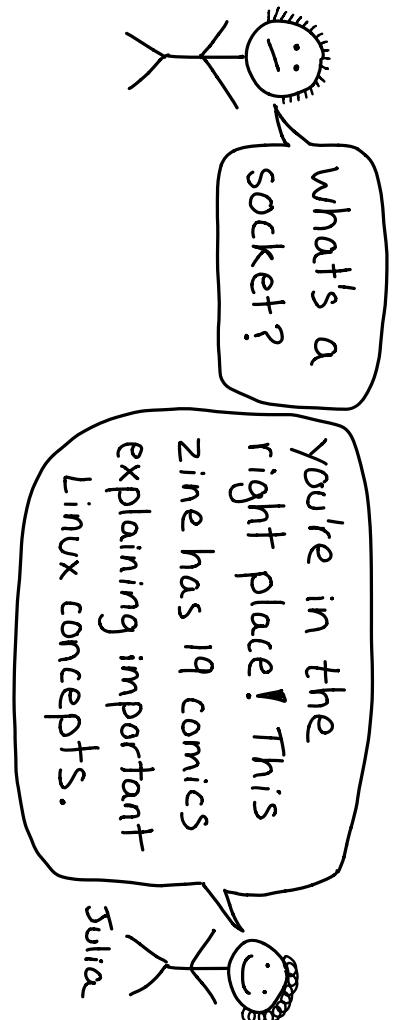
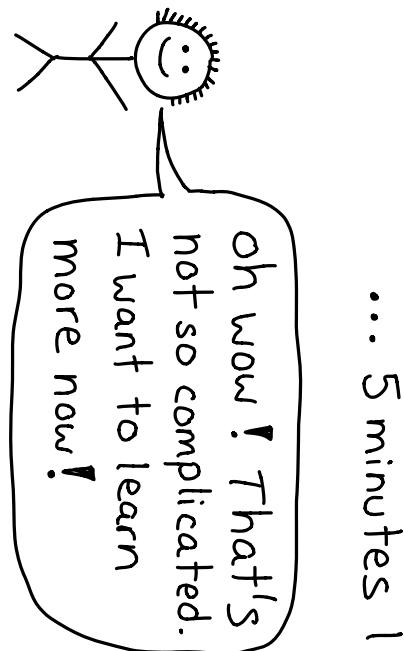


I want to learn more?
I highly recommend
this book: →

Every chapter is a readable,
short (usually 10-20 pages)
explanation of a Linux system.
I used it as a reference
constantly when writing
this zine.

I & it because even though
it's huge and comprehensive
(1500 pages!), the chapters
are short and self-contained
and it's very easy to pick it
up and learn something.

by Julia Evans
<https://jvns.ca>
twitter.com/b0rk



... 5 minutes later ...

man page sections 22

man pages are split up into 8 sections

① ② ③ ④ ⑤ ⑥ ⑦ ⑧

\$ man 2 read

means "get me the man page for **read** from section **2**".

There's both

→ a program called "read"

→ and a system call called "read"

so

\$ man 1 read

gives you a different man page from

\$ man 2 read

If you don't specify a section, man will look through all the sections & show the first one it finds.

man page sections

- ① programs \$ man grep \$ man ls
- ② system calls \$ man sendfile \$ man ptrace
- ③ C functions \$ man printf \$ man fopen
- ④ devices \$ man null for /dev/hull docs
- ⑤ file formats \$ man sudoers \$ man ls for /etc/sudoers \$ man sl
- ⑥ games \$ man proc files in /proc ! is my favourite from that section
- ⑦ miscellaneous \$ man apt explains concepts ! \$ man chroot
- ⑧ sysadmin programs \$ man symlink

Table of contents

unix permissions...4	sockets.....10	virtual memory...17
/proc5	unix domain sockets.....11	shared libraries...18
system calls.....6	processes.....12	copy on write19
signals.....7	threads.....13	page faults.....20
file descriptors.8	floating point....14	mmap.....21
pipes.....9	file buffering15	man page sections.....22
		memory allocation.....16

unix permissions

上

There are 3 things you can do to a file

read write execute

```

graph TD
    bork[bork]
    staff[staff]
    ANYONE[ANYONE]
    bork -- "rw--" --> bork
    staff -- "rw--" --> staff
    ANYONE -- "r---" --> ANYONE

```

The diagram illustrates file permissions for two groups: 'bork' and 'staff'. The 'bork' group has 'rw-' (read and write) permissions, while the 'staff' group also has 'rw-' permissions. The 'ANYONE' group has 'r--' (read only) permissions.

	setuid	setgid	user	group	all
000			110	110	100
sticky			r w x	r w x	r w x
For files:					
r	= can read				
w		= can write			
x			= can execute		
For directories, it's approximately:					
r	= can list files				
w		= can create files			
x			= can cd into & access files		

110 in binary is 6
So rw - r--
 110 100 100
 = 6 4
 = 4

chmod 644 file.txt
means change the
permissions to:
rw- r-- r--
Simple!

setuid affects executables

```
$ ls -l /bin/ping  
rws r-x r-x root root  
this means ping always runs as root
```

setgid does 3 different unrelated things for executables, directories, and regular files.

unix!
why???

it's a
long
story

UNIX

mm ap

21

What's mmap for?
I want to work with
a **VERY LARGE FILE**
but it won't fit
in memory.

You could
try mnmap!

With mmap
we all want to
read the same file!
no problem! mmap

Even if 10 processes mmap a file, it will only be read into memory once.

The diagram illustrates the relationship between a program and its components in memory. A large box labeled "Virtual memory" contains a smaller box labeled "RAM". Inside RAM, there are two sections: "Standard library" and "Dynamic linker". An arrow points from the "Standard library" section to a callout bubble containing the text "I need to use libc.so.6". Another arrow points from the "Dynamic linker" section to a callout bubble containing the text "I need to use ld". A third callout bubble, containing the text "you too eh? no problem.", has arrows pointing from both the "Standard library" and "Dynamic linker" sections.

Virtual memory

RAM

Standard library

Dynamic linker

I need to use libc.so.6

I need to use ld

you too eh? no problem.

load files lazily
with mmap

When you mmap a file, it
gets mapped into your
program's memory.

 ← 2TB of
virtual memory

but nothing is ACTUALLY
read into RAM until you
try to access the memory.
(how it works: page faults!)

<h2>how to mmap in Python</h2> <pre>import mmap f = open("HUGE.txt") mm = mmap.mmap(f.fileno(), 0)</pre> <p>↑ this won't read the file from disk! Finishes ~instantly.</p> <pre>print(mm[-1000:])</pre> <p>↑ this will read only the last 1000 bytes!</p> <h2>anonymous memory maps</h2> <ul style="list-style-type: none">→ not from a file (memory set to 0 by default)→ with MAP_SHARED, you can use them to share memory with a subprocess!	
--	--

page faults

20

every Linux process has a page table *

* page table *

virtual memory | physical memory address

0x19723000 0x1420000

0x19724000 0x1423000

0x1524000 not in memory

0x1844000 0x4a000 read only

some pages are marked as either

* read only

* not resident in memory

when you try to access a page that's marked "not resident in memory", it triggers a !page fault!

"not resident in memory" usually means the data is on disk !

virtual memory

in RAM

on disk

Having some virtual memory that is actually on disk is how swap and mmap work.

What happens during a page fault?

- the MMU sends an interrupt
- your program stops running
- Linux kernel code to handle the page fault runs

I'll fix the problem and let your program keep running

how swap works

① run out of RAM

RAM → disk

② Linux saves some RAM data to disk

RAM → disk

③ mark those pages as "not resident in memory"

in the page table not resident

virtual memory → RAM

④ when a program tries to access the memory, there's a !page fault!

⑤ time to move some data back to RAM!

Linux

virtual memory → RAM

RAM → virtual memory

⑥ if this happens a lot, your program gets VERY SLOW

I'm always waiting for data to be moved in & out of RAM

an amazing directory: /proc

/proc/PID/exe

symlink to the process's binary
magic : works even if the binary has been deleted!

/proc/PID/status

Is the program running or asleep? How much memory is it using? And much more!

/proc/PID/cmdline

command line arguments the process was started with

/proc/PID/environ

all of the process's environment variables

/proc/PID/fd

Directory with every file the process has open!

Run \$ ls -l /proc/42/fd to see the list of files for process 42.

These symlinks are also magic & you can use them to recover deleted files ♥

/proc/PID/stack

The kernel's current stack for the process. Useful if it's stuck in a system call.

/proc/PID/maps

List of process's memory maps. Shared libraries, heap, anonymous maps, etc.

for more information!

man proc

System calls

6

The Linux kernel has code to do a lot of things

`read from a hard drive`

`make network connections`

`create new process`

`change file permissions`

`kill process`

`change keyboard drivers`

every program uses system calls

`I use the 'open' syscall to open files`

`Python program`

`Java program`

`C program`

and every system call has a number (e.g. chmod is #90 on x86-64)
So what's actually going on when you change a file's permissions is:

`run Syscall #90 with these arguments`

`ok!`

`Linux`

TCP? dude I have no idea how that works.
No, I do not know how the ext4 filesystem is implemented. I just want to read some files!

your program doesn't know how to do those things

programs ask Linux to do work for them using `=system calls=`

`please write to this file`

(switch to running kernel code)

`done! I wrote 1097 bytes!`

`Linux`

`<program resumes>`

you can see which system calls a program is using with `=strace=`

`$ strace ls /tmp`

will show you every system call 'ls' uses!

it's really fun!

strace has high overhead

so don't run it on your production database

COPY ON WRITE

19

On Linux, you start new processes using the `fork()` or `clone()` system call.

calling fork creates a child process that's a copy of the caller

`Parent`

`child`

so Linux lets them share physical RAM and only copies the memory when one of them tries to `write`

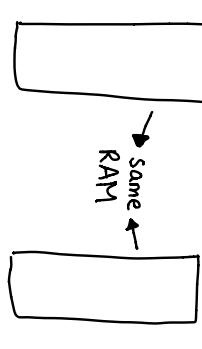
`I'd like to change that memory`

`Okay! I'll make you your own copy!`

`Linux`

the cloned process has EXACTLY the same memory.

- same heap
- same stack
- same memory maps if the parent has 3GB of memory, the child will too.



Linux does this by giving both the processes identical page tables.

copying all that memory every time we fork would be slow and a waste of RAM

often processes call `exec` right after `fork`, which means they don't use the parent process's memory basically at all!

When a process tries to write to a shared memory address:

- ① there's a `=page fault=`
- ② Linux makes a copy of the page & updates the page table
- ③ the process continues, blissfully ignorant

`It's just like I have my own copy`

shared libraries

18

Most programs on Linux use a bunch of C libraries:

`openssl` | `sqlite` (embedded db!)

`libpcre` | `zlib` (`gzip`!)

`libstdc++` (C++ standard library!)

how can I tell what shared libraries a program is using?

```
$ ldd /usr/bin/curl  
libz.so.1 => /lib/x86_64...  
libresolv.so.2 =>...  
libc.so.6 =>...  
+ 34 more !!
```

There are 2 ways to use any library:

① Link it into your binary

`your code` `zlib` `sqlite`

big binary with lots of things!

② Use separate shared libraries

`your code` ↗ all different `zlib` `sqlite` files

Programs like this

`your code` `zlib` `sqlite`

are called "statically linked"

and programs like this

`your code` `zlib` `sqlite`

are called "dynamically linked"

Where the dynamic linker looks

- ① DT_RPATH in your executable
- ② LD_LIBRARY_PATH
- ③ DT_RUNPATH in executable
- ④ /etc/ld.so.cache (run ldconfig -p to see contents)
- ⑤ /lib, /usr/lib

signals

7

You can send signals yourself with the `kill` system call or command

`SIGINT` `ctrl-C` } various levels of
`SIGTERM` `kill -9` } "die"
`SIGKILL` `kill -q` }
`SIGHUP` `kill -HUP` often interpreted as "reload config", e.g. by nginx

Signals can be hard to handle correctly since they can happen at ANY time

" oo handling a signal process
SURPRISE! Another signal!"

the Linux Kernel sends processes signals in lots of situations

that pipe is closed
your child terminated
illegal instruction
the timer you set expired
Segmentation fault

If you've ever used `Kill` you've used signals

DIE!!! Okay process

Your program can set custom handlers for almost any signal exceptions:

SIGSTOP & SIGKILL got → SIGKILLED
stop process
resume process

file descriptors

8

Unix systems use integers to track open files

`[!] Open foo.txt`

`[!] process`

`[!] kernel`

`[!] okay! That's file #7 for you.`

these integers are called **file descriptors**

When you read or write to a file / pipe / network connection you do that using a file descriptor

`[!] connect to google.com`

`[!] fd is 5!`

`[!] OS`

`[!] write GET / HTTP/1.1 to fd #5`

`[!] done!`

FD is for file descriptor

Let's see how some simple Python code works under the hood:

Python:

`f = open("file.txt")`

`f.readlines()`

Behind the scenes:

`[!] open file.txt`

`[!] Python program`

`[!] file #4`

`[!] here are the contents!`

`[!] OS`

`[!] ok! fd is 4`

`[!] read from stdin`

`[!] means`

"read from the file descriptor 0"

could be a pipe or file or terminal

`lsof (list open files) will show you a process's open files`

\$ lsof -p 4242 ← PID we're interested in

FD NAME

1 /dev/pts/4tuy1

2 pipe: 29174

3 /home/bark/awesome.txt

5 /tmp/ ↗

file descriptors can refer to:

→ files on disk

→ pipes

→ sockets (network connections)

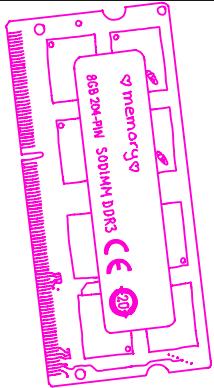
→ terminals (like xterm)

→ LOTS MORE (eventfd, inotify, signalfd, epoll, etc.)

Virtual memory

17

your computer has physical memory



Linux keeps a mapping from virtual memory pages to physical memory pages called the **page table**

a "page" is a 4kb or bigger chunk of memory

PID virtual addr physical addr
1971 0x20000 0x192000
2310 0x20000 0x228000
2310 0x21000 0x9788000

physical memory has addresses, like 0 - 8GB

but when your program references an address like 0x5c69a2a2, that's not a physical memory address!

It's a **virtual** address.

when your program accesses a virtual address

`[!] I'm accessing 0x21000`

CPU

`[!] I'll look that up in the page table and then access the right physical address`

MMU (Memory management unit)

hardware

every time you switch which process is running, Linux needs to switch the page table

`[!] program 1 0x129520 → "puppies"`

`[!] program 2 0x129520 → "bananas"`

`[!] Linux`

`[!] here's the address of process 2950's page table`

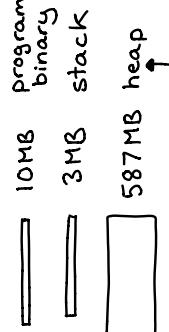
`[!] MMU`

`[!] thanks, I'll use that now!`

memory allocation

16

Your program has memory



Your memory allocator (malloc) is responsible for 2 things.

THING 1: Keep track of what memory is used/free.

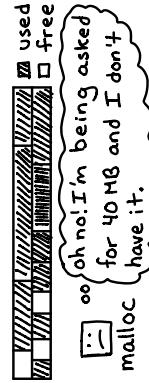


Your memory allocator's interface allocate `size_t size` bytes of memory & return a pointer to it.

`malloc (size_t size)`

mark the memory as unused (and maybe give back to the OS).
`realloc (void* pointer, size_t size)`
ask for more/less memory for pointer.
`calloc (size_t members, size_t size)`
allocate array + initialize to 0.

THING 2: Ask the OS for more memory!

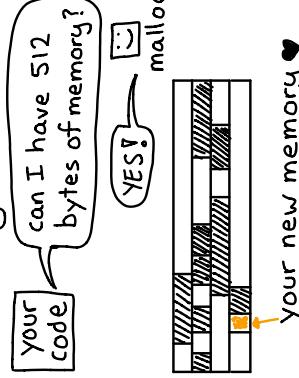


malloc isn't magic!!
it's just a function!

You can always:

- use a different malloc library like `jemalloc` or `tcmalloc` (easy!)
- implement your own malloc (harder)

malloc tries to fill in unused space when you ask for memory



pipes

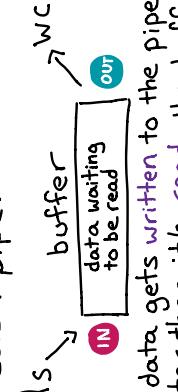
9

Sometimes you want to send the output of one process to the input of another

\$ ls | wc -l

53 ↪ 53 files!

Linux creates a buffer for each pipe.



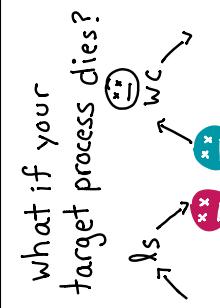
If data gets written to the pipe faster than it's read, the buffer will fill up. When the buffer is full, writes to `IN` will block (wait) until the reader reads. This is normal & ok.

when `ls` does
write(`IN`, "hi"),
`WC` can read it!
read(`OUT`)
→ "hi"

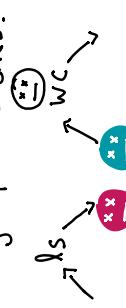
Pipes are one way. →
You can't write to `OUT`.

named pipes

\$ mkfifo my-pipe
This lets 2 unrelated processes communicate through a pipe!
f=open("./my-pipe")
f.write("hi!\n")
f=open("./my-pipe")
f.readline() ← "hi!"



what if your target process dies?

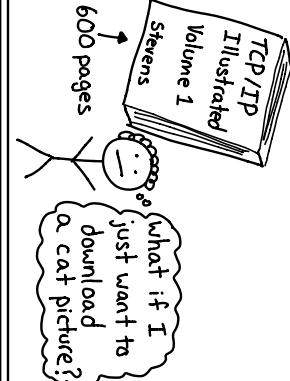


If `wc` dies, the pipe will close and `ls` will be sent `SIGPIPE`. By default, `SIGPIPE` terminates your process.

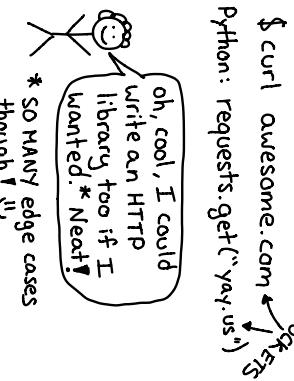
Sockets

10

networking protocols are complicated



Every HTTP library uses sockets under the hood



Unix systems have an API called the "socket API" that makes it easier to make network connections

- ① Create a socket
 - ② Connect to an IP/port
 - ③ Make a request
 - ④ Read the response
- ```
fd = socket(AF_INET, SOCK_STREAM)
connect(fd, ("12.3.4.5", 80))
write(fd, "GET /cat.png HTTP/1.1")
read(fd, ...)
```

AF-INET?  
What's that?

AF-INET means basically "internet socket": it lets you connect to other computers on the internet using their IP address.

The main alternative is AF-UNIX ("unix domain socket") for connecting to programs on the same computer.

SOCK\_STREAM = TCP  
curl uses this

SOCK\_DGRAM = UDP  
dig (DNS) uses this

SOCK\_RAW = just let me send IP packets.  
ping uses this  
I will implement my own protocol.

3 kinds of internet (AF-INET) sockets:

## file buffering

15

I printed some text but it didn't appear on the screen. why??

time to learn about flushing!

please write "I ❤ cats"  
process + file #1 (stdout)

• Write •

Linux

I/O libraries don't always call **write** when you **print**.

`printf("I ❤ cats");`

I'll wait for a newline before actually writing

This is called **buffering** and it helps save on syscalls.

- ① None. This is the default for `stderr`.
- ② Line buffering.  
(write after newline). The default for `terminals`.
- ③ "full" buffering.  
(write in big chunks)  
The default for `files` and `pipes`.

3 kinds of buffering (defaults vary by library)

default for `stderr`.

Line buffering.  
(write after newline).

default for `terminals`.

"full" buffering.  
(write in big chunks)

The default for `files` and `pipes`.

## flushing

Linux

To force your I/O library to write everything it has in its buffer right now, call **flush**!

studio

I'll call **write** right away!!

When it's useful to flush → when writing an interactive prompt!  
Python example:  
`print("password:", flush=True)`

→ when you're writing to a pipe/socket

no seriously, actually write to that pipe please

# floating point

14

a double is 64 bits  
 sign exponent fraction  
 $\downarrow$   $\uparrow$   
 100111 100111 100111 100111 100111 100111 100111  
 $E$ -1023  $\times$  1.frac  
 $\pm 2$

That means there are  $2^{64}$  doubles.  
 The biggest one is about  $2^{1023}$

weird double arithmetic

|                                                  |              |                                                        |
|--------------------------------------------------|--------------|--------------------------------------------------------|
| $2^{52} + 0.2 = 2^{52}$                          | $\leftarrow$ | (the next number after $2^{52}$ is $2^{52+1}$ )        |
| $1 + \frac{1}{52} = 1$                           | $\leftarrow$ | (the next number after $1$ is $1 + \frac{1}{2^{52}}$ ) |
| $2^{2000} = \text{infinity}$                     | $\leftarrow$ | infinity is a double                                   |
| $\text{infinity} - \text{infinity} = \text{nan}$ | $\leftarrow$ | nan = "not a number"                                   |

doubles get farther apart as they get bigger between  $2^n$  and  $2^{n+1}$  there are always  $2^{52}$  doubles, evenly spaced.  
 that means the next double after  $2^{60}$  is  $2^{60} + 64 = \frac{2^{60}}{2^{52}}$

they're very logical!  
 just understand how they work and don't use integers over  $2^{53}$  in Javascript ♥

they're scary and their arithmetic is weird!

they're very logical!  
 just understand how they work and don't use integers over  $2^{53}$  in Javascript ♥

they're very logical!  
 just understand how they work and don't use integers over  $2^{53}$  in Javascript ♥

## unix domain sockets !!

unix domain sockets are files.

\$ file mysock.sock  
 socket

the file's permissions determine who can send data to the socket.

they let 2 programs on the same computer communicate.

Docker uses Unix domain sockets, for example!

HTTP request  
 GET /container  
 Here you go! ~~~~~~

### advantage 1

Lets you use file permissions to restrict access to HTTP/ database services!

chmod 600 secret.sock

This is why Docker uses a unix domain socket.    

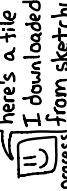
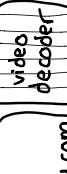
### advantage 2

UDP sockets aren't always reliable (even on the same computer).  
 unix domain datagram sockets are reliable!  
 And they won't reorder packets!

  
 and I know it'll arrive

### advantage 3

You can send a file descriptor over a unix domain socket.  
 Useful when handling untrusted input files!

  
  
 sandboxed process

# What's in a process?

12

|                                                                                                                                             |                                                                                                                           |                                                                                                          |                                                                                                                                |
|---------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------|
| <b>PID</b><br>process #129<br>reporting for duty!                                                                                           | <b>USER and GROUP</b><br>Who are you running as?<br>julia!                                                                | <b>ENVIRONMENT VARIABLES</b><br>like PATH! you can set them with:<br><code>\$ env A=val ./program</code> | <b>SIGNAL HANDLERS</b><br>I ignore SIGTERM!<br>I shut down safely!                                                             |
| <b>WORKING DIRECTORY</b><br>Relative paths ( <code>./blah</code> ) are relative to the working directory!<br><code>chdir</code> changes it. | <b>PARENT PID</b><br><code>pid 1 (init)</code> is everyone's ancestor<br><code>↑ PID 147</code><br><code>↑ PID 129</code> | <b>COMMAND LINE ARGUMENTS</b><br>See them in <code>/proc/pid/cmdline</code>                              | <b>OPEN FILES</b><br>Every open file has an offset.                                                                            |
| <b>MEMORY</b><br>heap! stack!<br>Shared libraries!<br>the program's binary!<br><code>mmaped</code> files!                                   | <b>THREADS</b><br>sometimes one<br>sometimes LOTS                                                                         | <b>CAPABILITIES</b><br>I have <code>CAP_PTRACE</code><br><code>Well I have CAP_SYS_ADMIN</code>          | <b>NAMESPACES</b><br><code>I'm in the host network namespace</code><br><code>I have my own namespace!</code> container process |
|                                                                                                                                             |                                                                                                                           |                                                                                                          |                                                                                                                                |

## threads

13

|                                                                                                                                                                                                                                                                              |                                                                                                                                                                                                                                                                   |                                                                                                                                                                                                                                                                                  |                                                                                                                                                                                                                                                                                                        |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p>Threads let a process do many different things at the same time</p> <p>process:</p> <ul style="list-style-type: none"> <li>thread 1: I'm calculating ten million digits of <math>\pi</math>! so fun!</li> <li>thread 2: I'm finding a REALLY BIG prime number!</li> </ul> | <p>threads in the same process share memory</p> <p>thread 1: I'll write some digits of <math>\pi</math> to 0x129420 in memory</p> <p>uh oh! that's where I was putting my prime numbers.</p> <p>thread 2: uh oh! that's where I was putting my prime numbers.</p> | <p>and they share code</p> <p><code>calculate-pi</code><br/> <code>find-big-prime-number</code></p> <p>but each thread has its own stack and they can be run by different CPUs at the same time</p> <p><code>[CPU 1] π thread</code><br/> <code>[CPU 2] primes thread</code></p> | <p>Why use threads instead of starting a new process? → a thread takes less time to create.</p> <p>→ sharing data between threads is very easy. But it's also easier to make mistakes with threads.</p> <p><code>[num] thread 1</code><br/> <code>You weren't supposed to CHANGE that data!</code></p> |
| <p>sharing memory can cause problems (race conditions!)</p> <p>memory</p> <p>thread 1: I'm going to add 1 to that number!</p> <p>thread 2: I'm going to add 1 to that number!</p> <p>RESULT: 24 ← WRONG. Should be 25!</p>                                                   | <p>thread 1: I'm going to add 1 to that number!</p> <p>thread 2: I'm going to add 1 to that number!</p>                                                                                                                                                           |                                                                                                                                                                                                                                                                                  |                                                                                                                                                                                                                                                                                                        |